

**DETERMINATION OF R-3-HYDROXYACYL-  
ACP-COA TRANSFERASE (PhaG) STRUCTURE**

**by**

**HASNI BIN ARSAD**

**Thesis submitted in fulfillment of the requirements  
for the degree of  
Doctor of Philosophy**

**March 2010**

## ACKNOWLEDGEMENT

First of all, I would like to thank my supervisor Prof. Madya Dr. Mohd Razip Samian and my co-supervisor Prof. Madya Dr. Habibah A. Wahab for their guidance and full support during the course of this study.

I would also like to express my sincere gratitude to Prof. Dr. Nazalan Najimudin and Prof. Maqsudul Alam (CCB@USM) for their kind support. To Dr. Irene Newhouse and especially, Dr. James Newhouse my heartfelt thanks for his patient tutoring in Molecular Dynamics, helpful discussion and support, particularly in molecular dynamics simulation and access to bioinformatic tools and the Maui supercomputer at the Universiti of Hawaii.

To Prof. Nor Muhammad Mahadi (MGI), Prof. Aishah A. Latiff (DCC-USM), Prof. Dr. Raja Noor Zaliha Raja Abd. Rahman (UPM), Dr Aida, Dr. Rashidah, Dr. Ahmad Sofiman Osman, Dr. Mustapha Fadzil, Dr. Tengku Sifzizul Tengku Mohamad, Dr. Amirul Ashraf and Dr. Alex; my heartfelt thanks for their advice and permission to use their chemicals and lab apparatus.

I am grateful to the Ministry of Science, Technology and Innovation, Malaysia for providing me the National Science Fellowship. I am also grateful to MGI for support in terms of scientific activities and access to their supercomputer.

To members of Lab 414, Lab 406, Lab 409, Lab 318 and CCB; thank you for your direct or indirect contribution toward the accomplishment of my PhD study.

Last but not least, my heartfelt gratitude to Salina Mokhtar, Ameerah Husna Syafiqah, Ameer Syahir Hafidzi, Ameerah Hanim Syuhada and my parents (Arsad Hasan and Samsiah) whose continuous and unwavering support permit me to pursue my ambitions.

## TABLE OF CONTENTS

ACKNOWLEDGEMENT .....	ii
TABLE OF CONTENTS .....	iii
LIST OF TABLES .....	vi
LIST OF FIGURES .....	viii
LIST OF SYMBOLS, ABBREVIATIONS AND NOMENCLATURES .....	xii
ABSTRAK .....	xv
ABSTRACT .....	xvii
CHAPTER 1 .....	1
1.0 INTRODUCTION .....	1
1.1 Research Objectives.....	3
1.2 Thesis overview .....	3
1.3 Literature review .....	4
1.3.1 Protein structure: .....	4
1.3.2 Protein structure determination .....	20
1.3.3 Protein structure database.....	21
1.3.4 Protein structure prediction .....	22
1.3.5 Protein structure validation: .....	28
1.3.6 Enzymes .....	30
1.3.7 The enzyme catalysis.....	33
1.3.8 PHA biosynthesis in <i>Pseudomonas</i> .....	33
1.3.9 Fatty acid <i>de novo</i> synthesis .....	37
1.3.10 (R)-3-hydroxyacyl-ACP – CoA tranferase (PhaG) enzyme.....	41
1.4 PHA in transgenic plants .....	54
1.4.1 Production of PHAs in transgenic plants.....	54

CHAPTER 2.....	57
EXPRESSION AND PURIFICATION OF R-3-HYDROXYACYL-ACYL CARRIER PROTEIN COENZYME A TRANSFERASE (PhaG) OF <i>Pseudomonas</i> <i>sp.</i> USM 4-55	
2.0 INTRODUCTION.....	57
2.1 OBJECTIVES OF THIS STUDY.....	61
2.2 METHOD .....	61
2.2.1 Bacterial strain and plasmid .....	61
2.2.2 General methods.....	63
2.2.3 Culture medium.....	64
2.2.4 Cloning <i>phaG</i> <sub><i>P. sp.</i> USM 4-55</sub> into the expression vector .....	64
2.2.5 Expression Vector Preparation.....	66
2.2.6 Ligation of PCR product into expression vector .....	66
2.2.7 Transformation .....	66
2.2.8 Confirmation of expression construction by DNA sequencing.....	71
2.2.9 Analysis of nucleotide sequence .....	71
2.2.10 Protein expression and purification.....	71
2.2.11 SDS-PAGE analysis .....	72
2.2.12 Protein Concentration determination.....	73
2.2.13 Determination of PhaG <sub><i>P. sp.</i> USM 4-55</sub> Isoelectric Point (pI) value.....	73
2.2.14 Silver Staining .....	76
2.2.15 Mass Spectrometry Analysis of PhaG <sub><i>P. sp.</i> USM 4-55</sub> .....	79
2.2.16 Protein crystallization.....	85
2.3 RESULTS AND DISCUSSION.....	87
2.3.1 Overexpression of PhaG in recombinant <i>Escherichia coli</i> .....	87
2.3.2 Expression of PhaG in <i>E. coli</i> .....	89
2.3.3 Protein purification.....	93
2.3.4 Isoelectric Point (pI) of PhaG .....	97

2.3.5	Mass Spectrometry Analysis of trypsin treated of PhaG <sub>P. sp USM 4-55</sub> ...	97
2.3.6	Protein crystallization screening.....	116
CHAPTER 3:.....		118
PREDICTION OF PhaG STRUCTURE		
3.0	INTRODUCTION.....	118
3.1	OBJECTIVES OF THIS STUDY.....	123
3.2	METHODS .....	123
3.2.1	Prediction of PhaG <sub>P. sp USM 4-55</sub> structure .....	125
3.2.2	Secondary structure prediction by Circular Dichroism .....	128
3.2.3	Prediction of tertiary structure of PhaG <sub>P. sp USM 4-55</sub> by threading method .....	128
3.2.4	Docking simulation of the substrate on the PhaG <sub>P. sp USM 4-55</sub> .....	130
3.3	RESULT AND DISCUSSION .....	131
3.3.1	Sequence analysis.....	132
3.3.2	Secondary structure characterisation of PhaG <sub>P. sp USM 4-55</sub> .....	145
3.3.3	Template identification using BioInfobank Metaserver and Threader3.5.....	146
3.3.4	Model building and refinement .....	156
3.3.5	Validation of the PhaG <sub>P. sp USM 4-55</sub> model.....	162
3.3.6	Binding site prediction by SiteMap 2.2 (Schrödinger suit 2008).....	172
3.3.7	Docking Analysis using Glide (Schrödinger suite 2008) .....	175
3.3.8	Proposed mechanism of R-3-hydroxyacyl-ACP-CoA transferase ....	184
CHAPTER 4.....		190
4.0	GENERAL DISCUSSION .....	190
5.0	CONCLUSION .....	206
REFERENCES.....		209
APPENDICES.....		233

## LIST OF TABLES

Table 1.1. Standard amino acid abbreviation and side chain properties .....	10
Table 1.2. Relative amino acid compositions of mesophiles and thermophiles One-letter abbreviations of amino-acid residues in brackets .....	15
Table 2.1. Bacterial strains and plasmids .....	62
Table 2.2. Primer sequence for the cloning of <i>phaG<sub>P. sp. USM 4-55</sub></i> into the expression vector pQE-30 Restriction adapter sequences are colored red. ....	65
Table 2.3. PCR cycle conditions for the amplification of <i>phaG<sub>P. sp. USM 4-55</sub></i> .....	67
Table 2.4. Composition of PCR product digestion mix. ....	68
Table 2.5. Digestion of Cloning Vector (PQE 30, Qiagen).....	69
Table 2.6. Composition of ligation mixture between expression vector pQE 30 and insert (PCR product) .....	70
Table 2.7. Composition of SDS-PAGE gels .....	74
Table 2.8. IEF focusing condition as recommended by ReadyStrip™ IPG strip instruction manual (Bio-Rad). ....	77
Table 2.9. Composition of solutions used in silver staining of the polyacrylamide gel .....	78
Table 2.10. Parameters of LC/MS Q-TOF (Agilent 6250) used for mass spectrometry analysis of PhaG <sub>P. sp. USM 4-55</sub> .....	82
Table 2.11. Monoisotopic masses and immonium ions used to identify amino acid sequences from MS/MS spectra.....	83
Table 2.12. Basic rules for amino acid sequence determination .....	84
Table 2.13. Basic rules to determine a fragment series .....	84
Table 2.14. Purification table for PhaG <sub>P. sp. USM 4-55</sub> .....	98
Table 2.15. Peptide masses of PhaG treated with trypsin as predicted by ExPASy server ( <a href="http://www.expasy.ch/cgi-bin/peptide-mass.pl">http://www.expasy.ch/cgi-bin/peptide-mass.pl</a> ) .....	101
Table 2.16. The <i>de novo</i> sequencing by PEAKS studio 5.0 of the peptides derived from PhaG digested with trypsin. Fifteen peptides were identified, which comprised of 199 amino acid residues. The peptides were sorted sequentially according to the amino acid sequence of PhaG <sub>P. sp. USM 4-55</sub> . ....	105

Table 2.17. The manual <i>de novo</i> sequencing of the peptide of PhaG digested with trypsin. Four peptides, which comprised of 35 amino acids were identified. The peptides were sorted sequentially according to PhaG <sub>P. sp. USM 4-55</sub> amino acid sequence. ....	113
Table 2.18. The MS/MS spectra analysis by Mascot search server, Peaks studio 5.0 and manual <i>de novo</i> sequencing. Twenty-four amino acids that were not detected by the MS/MS spectra are shown against a red background.....	115
Table 3.1. Amino acid composition of PhaG <sub>P. sp. USM 4-55</sub> . (EXPASY Proteomic Server, <a href="http://www.expasy.ch/">http://www.expasy.ch/</a> ). ....	135
Table 3.2. Result of PhaG <sub>P. sp. USM 4-55</sub> BLASTP analysis using experimentally determined protein structure database at RCSB protein data bank....	139
Table 3.3. The species and the PhaG accession numbers found in GenBank .....	144
Table 3.4 The predicted secondary structure of PhaG <sub>P. sp. USM 4-55</sub> obtained from secondary structure servers PHD, PROF, PSIPRED, JNET, SSPPRO, PREDATOR and JUFO. C - indicates random coils, E - indicates extended or $\beta$ -sheets, H- indicates $\alpha$ -helix. CONS – indicates consensus secondary structure. ....	147
Table 3.5. The percentage of secondary structure composition of PhaG amino acid sequence predicted using secondary structure prediction server and CD spectroscopy.....	149
Table 3.6 Result of Threader 3.5 template search. Column 8 is the Z-score. ....	152
Table 3.7. Template search by BioInfoBank Metaserver ( <a href="http://meta.bioinfo.pl">http://meta.bioinfo.pl</a> ). .	153
Table 3.8. Secondary structure analysis of PhaG <sub>P. sp. USM 4-55</sub> . ....	166
Table 4.1. Validation of PhaG model 132 PhaG <i>P. putida</i> -1CQZ , PhaG <i>P. putida</i> -1EHY and PhaG <i>P. mendocina</i> -1EHY. ....	201

## LIST OF FIGURES

Figure 1.1. Basic structure of an amino acid. ....	5
Figure 1.2. Geometry of the peptide backbone. ....	7
Figure 1.3. Ramachandran plot angles. ....	8
Figure 1.4 The Ramachandran plot showing the dihedral angle at N-C $\alpha$ -C, i.e. the $\phi$ and $\psi$ angle for the peptide backbone in the protein structure. ....	9
Figure 1.5. A Venn diagram illustrating the properties of amino acids (Taylor, 1986, Betts and Russell, 2003) .....	11
Figure 1.6. A right-handed alpha helix. ....	16
Figure 1.7. Two forms of $\beta$ -sheets are commonly found in secondary structure, namely antiparallel and parallel $\beta$ -sheets. ....	17
Figure 1.8. Metabolic pathways involved in mcl-PHA precursor synthesis in <i>Pseudomonas</i> . ....	34
Figure 1.9. Pathway and components of fatty acid $\beta$ -oxidation identified by Pathway Studio programme analyzing the Ariadne Bacteria 2.01 database. ....	36
Figure 1.10. Metabolic pathway of fatty acid $\beta$ -oxidation in <i>Pseudomonas putida</i> KT2440 (KEGG). ....	38
Figure 1.11. Fatty acid $\beta$ -oxidation pathways that supply monomers to the mcl-PHA synthase in pseudomonads (Fiedler <i>et al.</i> , 2002). ....	39
Figure 1.12. Pathway and components of fatty acid biosynthesis identified by Pathway Studio programme analyzing the Ariadne Bacteria 2.01 database. ....	40
Figure 1.13. Fatty acid biosynthesis of <i>Pseudomonas putida</i> KT2440 (KEGG). ....	42
Figure 1.14. <i>De novo</i> fatty acid biosynthesis pathway that supplies monomers to PHA synthase (PhaC) (Fiedler <i>et al.</i> , 2000). ....	43
Figure 1.15. The PhaG pathway diagram created by Pathway Studio which translated the data collected by MedScan. ....	45
Figure 1.16. PHA synthase precursors from fatty acid biosynthesis metabolic pathway .....	46
Figure 1.17. The proposed model of PHA production in recombinant <i>E. coli</i> from sugars <i>via de novo</i> fatty acid biosynthesis. ....	50



Figure 1.18. The proposed pathways for PHA monomer supply from fatty acid biosynthesis in <i>E. coli</i> strains overproducing FabH (Nomura <i>et al.</i> , 2004). .....	52
Figure 1.19. Metabolic engineering of the fatty acid biosynthesis pathway for production of PHA (Park <i>et al.</i> , 2005). .....	53
Figure 1.20 Pathways for PHA formation in plant plastids. Solid arrows indicate native plant enzyme activities. ....	56
Figure 2.1. Amplification of <i>phaG</i> using primers PHAGH_Fxa and PHAGB-s. ....	88
Figure 2.2. Purified digestions ( <i>Bam</i> H1 and <i>Hind</i> III) of vector (pQE-30) and insert ( <i>phaG</i> <sub>P. sp USM 4-55</sub> ). ....	90
Figure 2.3. Plasmid map of pQHG-5.....	91
Figure 2.4. The nucleotide sequence of pQHG-5. ....	92
Figure 2.5. PCR of colonies suspected to carry the full-length <i>phaG</i> . ....	94
Figure 2.6. Protein expression time profile in colonies 3, 5 and 6. ....	95
Figure 2.7. SDS-PAGE analysis of purified PhaG <sub>P. sp. USM 4-55</sub> . ....	96
Figure 2.8. Silver stained IEF 2-D gel. ....	99
Figure 2.9. The matched peptides are shown in bold red. ....	103
Figure 2.10 Predicted peptide sequences of PhaG digested with trypsin. ....	104
Figure 2.11. The MS/MS spectrum of the peptide (M+H) <sup>+</sup> = 329.70120.....	107
Figure 2.12. This MS/MS spectrum showing y-ion series for peptide (M+H) <sup>+</sup> = 657.36972. ....	109
Figure 2.13. This MS/MS spectrum shows b-ion series for peptide (M+H) <sup>+</sup> = 657.36972. The glutamic acid (E) also appeared on b-ion series. ....	109
Figure 2.14. The MS/MS spectrum shows the b3 from b-ion series. The a ion (270 – 28=242) and c ion (270 + 17 = 287) is shown as well.....	111
Figure 2.15. The MS/MS spectrum shows the y4 from y-ion series. The z ion (490 – 17=473) and x ion (490 + 28 = 518) was identified as well. ....	111
Figure 2.16. This MS/MS spectrum show the immonium ions of the peptide (M+H) <sup>+</sup> 659.4024.....	112
Figure 2.17. The result of the protein crystallization screening. ....	117
Figure 3.1. A flowchart of PhaG structure prediction.....	124
Figure 3.2. The amino acid sequence of PhaG <sub>P. sp USM 4-55</sub> (295 residues). ....	133

Figure 3.3. The top 20 proteins that showed high homology to PhaG <sub>P. sp USM 4-55</sub> as a result of BLASTP analysis of PhaG <sub>P. sp USM 4-55</sub> on non-redundant protein database. ....	136
Figure 3.4. The alignment of PhaG to the $\alpha/\beta$ hydrolase fold protein family. ....	137
Figure 3.5 Alignment of the 21 PhaG amino acid sequences. Amino acids in red indicates conserved sequences. ....	140
Figure 3.6. An alignment between PhaG <sub>P. sp USM 4-55</sub> and its modeling template 1C4X using ClustalW. ....	155
Figure 3.7. Input file for homology modeling needed by Modeller9v4 was the sequence alignment between PhaG and 1C4X in PIR format. ....	158
Figure 3.8. Ramachandran plot from PROCHECK analysis of the model 42 as produced by Modeller9v4. ....	160
Figure 3.9. Ramachandran plot of model PhaG 132, after energy minimization process. ....	161
Figure 3.10. A cartoon of PhaG model 132 drawn using VMD (Humphrey <i>et al.</i> , 1996) visual software. ....	163
Figure 3.11. Superimposed structure of PhaG model 132 and 1C4X by Swiss-Pdb Viewer (Guex and Peitsch, 1997). ....	164
Figure 3.12 The Secondary structure analysis of PhaG model 132 using VMD 1.8.6. ....	165
Figure 3.13. A. The hydrophobic region of PhaG model 132 (blue). B. The polar residues (green), mostly covering the hydrophobic residues. ....	169
Figure 3.14. PhaG Model 132 after the loop refinement process using Modeller9v4. ....	171
Figure 3.15. The ligand binding site of PhaG <sub>P. sp USM 4-55</sub> as predicted by SiteMap. ....	173
Figure 3.16. The location of active site with respect to the ligand binding site. ....	174
Figure 3.17. The grid box was prepared with a size of 20x20x20 angstrom in length which covered the predicted active site and the binding pocket predicted by SiteMap. ....	177
Figure 3.18. The ligand, R-3-hydroxyoctanoate-CoA, used in the docking simulation. ....	178
Figure 3.19. The ligand docked into the predicted binding site. ....	179

Figure 3.20. Docking simulation showing the ligand docked at the predicted binding site and the positions of the catalytic triad residues (Ser-102, Asp-223 and His-251).....	180
Figure 3.21. The position of the ligand in the binding pocket. ....	182
Figure 3.22. The position and distance between the predicted catalytic triad residue Ser-102, His-251 and Asp-223, and the thioester bond of the ligand. .	183
Figure 3.23. The catalytic triad residues Ser-102, Asp-223 and His-251 of PhaG model 132 takes up similar positions when compared to the other catalytic triads of other serine proteinases and lipases. .	187
Figure 3.24. Conformations of representative catalytic triads from each of the four fold groups.....	187
Figure 3.25. PhaG catalytic mechanism. ....	189
Figure 4.1. Secondary structure diagram of PhaG model 132. ....	194
Figure 4.2. The amino acid residues (Trp-103, Met-135, Gly-142, Leu-157, Met-189, Val-196, Tyr-225) of the PhaG model which are involved in the hydrophobic contact with the ligand as analysed by LIGPLOT (Wallace <i>et al.</i> , 1995).....	197
Figure 4.3. Superimposition of PhaG model 132 over the model proposed by Zheng <i>et al.</i> (2005) and Hoffmann <i>et al.</i> (2002). ....	198
Figure 4.4. Superimposition of catalytic triads (Ser-His-Asp) acquired from the proposed model of PhaG model 132 (blue), PhaG <i>P. putida</i> -1CQZ (yellow), PhaG <i>P. putida</i> -1EHY (green) and PhaG <i>P. mendocina</i> -1EHY (red).....	200
Figure 4.5. The fatty acid elongation pathway in plants (www.plantcyc.org). ....	204

## LIST OF SYMBOLS, ABBREVIATIONS AND NOMENCLATURES

Å	Angstrom
$\alpha$	Alpha
$\beta$	Beta
°C	Degree Celsius
%	Percentage
$\phi$	Phi
$\psi$	Psi
3-D	Three dimensional
3HB-co-3HV	3-hydroxybutyrate-co-3-hydroxyvalerate
3HB	3-hydroxybutyrate
3HA	3-hydroxyacyl
ACP	Acyl carrier protein
ATP	Adenosine triphosphate
APS	Ammonium persulfate
CoA	Coenzyme A
DNA	Deoxyribonucleic acid
g	Gram
IEF	Isoelectric focusing
IPTG	Isopropyl $\beta$ -D-thiogalactopyranoside
LJ	Lennard Jones

lcl-PHA	Long chain length-PHA
LB	Luria Bertani
MCL	Medium chain length
mcl-PHA	Medium chain length polyhydroxyalkanoate
mcl-3HA	Medium chain length 3-hydroxyalkanoates
μm	Micrometer
μl	Microliter
μM	Micromole
μg/ml	Microgram per milliliter
mg/l	Milligram per liter
mM	Millimole
ml	Milliliter
M	Molar
MD	Molecular dynamics
nm	Nanometer
NADPH	Nicotinamide adenine dinucleotide phosphate
TEMED	N,N,N',N'-tetramethylethylenediamine
N	Normality
NMR	Nuclear magnetic resonance
OD	Optical density
ps	Picosecond
PCR	Polymerase chain reaction
PHA	Polyhydroxyalkanoates

PDB	Protein Data Bank
psi	Pounds per square inch
RMSD	Root mean square deviation
rpm	Rotation per minute
scl-PHA	Short chain length-PHA
SDS	Sodium dodecyl sulfate
SDS-PAGE	Sodium dodecyl sulfate polyacrylamide gel electrophoresis
TAE	Tris-acetate/EDTA electrophoresis buffer
UV	Ultra violet
U	Units
V	Volt
V-hr	Volt per hour

## PENENTUAN STRUKTUR R-3-HIDROKSIASIL-ACP-COA TRANSFERASE (PhaG)

### ABSTRAK

R-3-hidroksiasil-ACP-CoA transferase (PhaG) berfungsi sebagai pemangkin penukaran (R)-3-hidroksiasil-ACP kepada (R)-3-hidroksiasil-CoA, sebagai prekursor utama untuk pempolimeran polyhydroxyalkanoate (PHA) daripada substrat tidak berkaitan. PHA adalah bioplastik yang berpotensi untuk menggantikan termoplastik berasaskan fosil kerana ianya boleh di biourai. Enzim PhaG daripada pencilan strain tempatan, *Pseudomonas* sp. USM 4-55, telah berjaya diklon (nombor akses Genbank EU305558). Pada masa ini struktur 3 dimensi yang sama dengan PhaG masih tidak diketahui. Untuk ekspresi lampau gen *phaG* telah diklon ke dalam vector ekspresi, pQE ke-30 dan telah berjaya diekspresi melalui aruhan menggunakan 0.5 mM IPTG di dalam sel perumah *Escherichia coli* SG 13009. PhaG berjaya dimurnikan dengan menggunakan kaedah kromatografi afiniti logam diikuti oleh kromatografi penurasan gel dan analisis Coomassie Blue SDS-PAGE menunjukkan satu jalur. Protein yang telah dimurnikan itu telah dipastikan sebagai PhaG melalui kaedah penjujukan peptide *de novo*. Nilai pI bagi PhaG<sub>P. sp USM 4-55</sub> yang ditentukan oleh IEF adalah 7.45. Penskrinan hablur protein telah dilakukan menggunakan kaedah “microbatch” dan “hanging drop”. Malangnya, tidak satu pun daripada kit penghabluran yang digunakan menunjukkan hablur PhaG. Satu model 3-D PhaG dicadangkan dengan kaedah “threading”. Protin struktur enzim 2-Hydroxy-6-okso-6-phenylhexa-2,4 hydrolase daripada *Rhodococcus* sp. RHA1 (identity PDB: 1C4X) digunakan sebagai template. Modeller9v4 digunakan untuk menghasilkan model PhaG<sub>P. sp USM 4-55</sub>. Nilai RMS diantara model PhaG model 132 dan template, 1C4X adalah 0.75Å. Menariknya, analisis tapak pengikatan pada

PhaG model 132 yang telah ditentukan menggunakan SiteMap telah berjaya melabuhkan ligan R-3-hydroxyoctanoate-CoA di tapak pengikatan dengan Glide. Seterusnya mekanisme tindakbalas enzim PhaG dicadangkan berdasarkan analisis docking, dan mekanisme enzim malonil CoA-ACP transakilase (FabD).



## DETERMINATION OF R-3-HYDROXYACYL-ACP-COA TRANSFERASE (PhaG) STRUCTURE

### ABSTRACT

R-3-hydroxyacyl-ACP-CoA transferase (PhaG) catalyzes the conversion of (R)-3-hydroxyacyl-ACP to (R)-3-hydroxyacyl-CoA derivatives, which serve as the ultimate precursor for polyhydroxyalkanoate (PHA) polymerization from unrelated substrates. PHA is a family of bioplastic that has a good potential to replace fossil-based thermoplastics because it is biodegradable. The transferase enzyme PhaG of a locally isolated strain, *Pseudomonas* sp. USM 4-55, was recently cloned (GeneBank accession number EU305558). Currently there is no known 3D structure with high similarity to PhaG. In order to over express, the phaG gene was cloned into expression vector pQE-30 and it was successfully overexpressed by induction with 0.5 mM IPTG in the host *Escherichia coli* strain SG 13009. The PhaG was purified by metal affinity chromatography followed by gel filtration chromatography to a single band under Coomassie Blue SDS-PAGE analysis. The purified protein was confirmed as PhaG by *de novo* peptide sequencing. The pI of PhaG<sub>P. sp USM 4-55</sub> is 7.45, as determined by IEF. Protein crystallization screening was then performed by applying microbatch and hanging drop method. However, none of the crystallization screening kit used were able to produce PhaG protein crystal. A 3-D model of PhaG was predicted using threading method. The 2-hydroxy-6-oxo-6-phenylhexa-2,4-dienote hydrolase enzyme from *Rhodococcus* sp. strain RHA1 (pdb id: 1C4X) structure was used as the template. The Modeller9v4 program was used to model the PhaG<sub>P. sp USM 4-55</sub>. The RMS value of the PhaG model 132, and the template 1C4X is 0.75Å. Interestingly, binding site prediction using SiteMap on the proposed model, showed a binding pocket that can accomodate the ligand, R-3-hydroxyoctanoate-

CoA by using Glide. Consequently, the reaction mechanism of PhaG is proposed based on the docking analysis and the mechanism of malonyl CoA-ACP transacylase (FabD).

## CHAPTER 1

### 1.0 INTRODUCTION

Polyhydroxyalkanoate (PHA) is a kind of bioplastic that has a good potential to replace fossil-based thermoplastics due to its biodegradable properties. PHAs are normally found in various microorganisms, accumulated as intracellular granules, that can serve as carbon and energy storage material under nutrient-limiting conditions (Dawes and Senior, 1973).

PHA can be divided into three groups depending on the number of carbon atoms in the monomer units. The first two groups are short-chain-length PHA (scl-PHA), which consist of 3 – 5 carbon atoms and medium-chain-length PHA (mcl-PHA), which consist of 6 – 14 carbon atoms. The third group, long-chain-length PHA (lcl-PHA), is reserved for polymers containing more than 14 carbon atoms (Steinbüchel *et al.*, 1992).

Steinbüchel (2001) reported that there are approximately 150 different hydroxyalkanoic acids known to be constituents of bacterial storage polyester (PHA). Bacterial polyesters such as the homopolyester poly(3HB), the copolyester poly(3HB-co-3HV) and PHAs consisting of 3-hydroxyoctanoate, 3-hydroxydecanoate and a few other medium-chain-length 3-hydroxyalkanoates (poly(mcl-3HA)) have been manufactured and used in various applications. Byrom (1994) reported the commercial production of poly(3HB) and poly(3HB-co-3HV) as biodegradable bioplastics. PHAs are also used in product applications such as latex paints (van der Walle *et al.*, 2001), medical application such as retard material (Fraser *et al.*, 1989) and as scaffolding material for tissue engineering (Williams *et al.*, 1999).

The unique and interesting properties of PHAs attracted academic and commercial bodies to embark on extensive research to study this biopolymer. Biosynthesis of PHA, especially in the bacterial cell was investigated by a few groups of scientist. A diverse range of bacteria are able to accumulate PHA, for example almost all *Pseudomonads* synthesize poly(mcl-3HA) when cultured on alkanes, organic acids, glucose and many other carbon sources (Haywood *et al.*, 1990, Timm and Steinbuchel, 1990).

Huijberts *et al.* (1994) and Rehm *et al.* (1998) found that there are at least three different metabolic routes in *P. putida* for the synthesis of 3-hydroxyacyl coenzyme A, the substrate for PHA synthase to synthesize PHA. They are 1) beta-oxidation pathway, 2) fatty acid *de novo* biosynthesis pathway and 3) chain elongation reaction pathway.

A further investigation on PHA synthesis by fatty acid *de novo* revealed that R-3-hydroxyacyl-acyl carrier protein-Coenzyme A transferase (PhaG) was the enzyme responsible to channel down substrates from fatty acid *de novo* biosynthesis pathway to the PHA synthase for PHA accumulation in *P. putida* (Rehm *et al.*, 1998). The evidence showed that PhaG catalyzes the conversion of (R)-3-hydroxyacyl-ACP to (R)-3-hydroxyacyl-CoA derivatives which serve as the ultimate precursor for PHA polymerization from unrelated substrates.

Van der Leij and Witholt (1995) were the first to allude to the connection between *de novo* fatty acid biosynthesis and PHA production in plant. They predicted the existence of a transferase enzyme that was needed to convert (R)-hydroxyacyl-ACP available in the plant fatty acid biosynthesis pathway to the (R)-hydroxyacyl-CoA, which is the substrate for PHA synthase.

In this thesis, the PhaG from *Pseudomonas* sp. USM 4-55 (PhaG<sub>P. sp. USM 4-55</sub>) was used as the subject. This strain was isolated from soil samples taken from Tasek Chini, Pahang (Few, 2001). This isolate was chosen for this thesis because of its rare characteristic; it was able to accumulate two groups of polymer at the same time, namely scl-PHA and mcl-PHA. Previous studies showed that *P. sp* USM 4-55 could accumulate mcl-PHA up to 20 % of cell dry weight when cultured on 2 % glucose (Few, 2001). The result of the study mentioned above indicates that the PhaG in this strain is functional as described by Rehm *et al.*, (1998). PhaG plays an important role in channeling the substrates for PHA synthase when glucose or gluconate is used as the carbon source (Rehm *et al.*, 1998). This enzyme is usually present in *Pseudomonads*. The unique and interesting characteristic of PhaG attracted me to embark on an effort to determine its structure and understand the functional mechanism of this enzyme.

## **1.1 Research Objectives**

1. To overexpress and crystallise PhaG
2. To model 3-dimensional structure of PhaG
3. To propose a reaction mechanism of PhaG

## **1.2 Thesis overview**

The thesis is divided into four chapters. It starts of with a literature review of protein structure in general and PhaG in the chapter 1. The second chapter involves the overexpression of PhaG<sub>P. sp USM 4-55</sub> in *E. coli*. Details about cloning, expression and purification, pI determination and protein crystal screening will be elaborated.

Chapter three is an analysis to understand the three dimensional structure of this protein by employing bioinformatics tools to build and predict the protein model. In addition, the reaction mechanism of PhaG is also proposed. Chapter four is a general discussion and the suggestion for the future studies.

### **1.3 Literature review**

#### **1.3.1 Protein structure:**

Protein has biologically diverse functions in the cell ranging from DNA replication, forming cytoskeletal structures, transporting oxygen around the bodies of multicellular organism to converting one molecule into another. Lately, the studies of protein structures and functions have increased. This is due to the current trend to understand of protein interactions with other biomolecules and their role within different biological systems and their potential manipulations by employing genetic or chemical methods.

Biochemists have classified protein structures into four different classes, namely primary structure, secondary structure, tertiary structure and quaternary structure. The primary protein structure is the amino acid sequence in a linear heteropolymer amino acid chain. Amino acid is consisted of atoms C, O, N, H and S. The basic structure of amino acid is shown in Figure 1.1. A carbon atom ( $C_{\alpha}$ ) is bound to an amino group ( $-NH_2$ ), a carboxyl group ( $-COOH$ ), a hydrogen atom (H) and an organic side group (assigned as R). A statistical survey of x-ray structures in the Cambridge Structure Database found that the average of C-N peptide bond length is 1.33 Å, except for proline residue, which is 1.34 Å; the C-O bond length is 1.23 Å,  $C_{\alpha}$ -C bond length is 1.53 Å except glycine residue, which is and 1.52 Å,

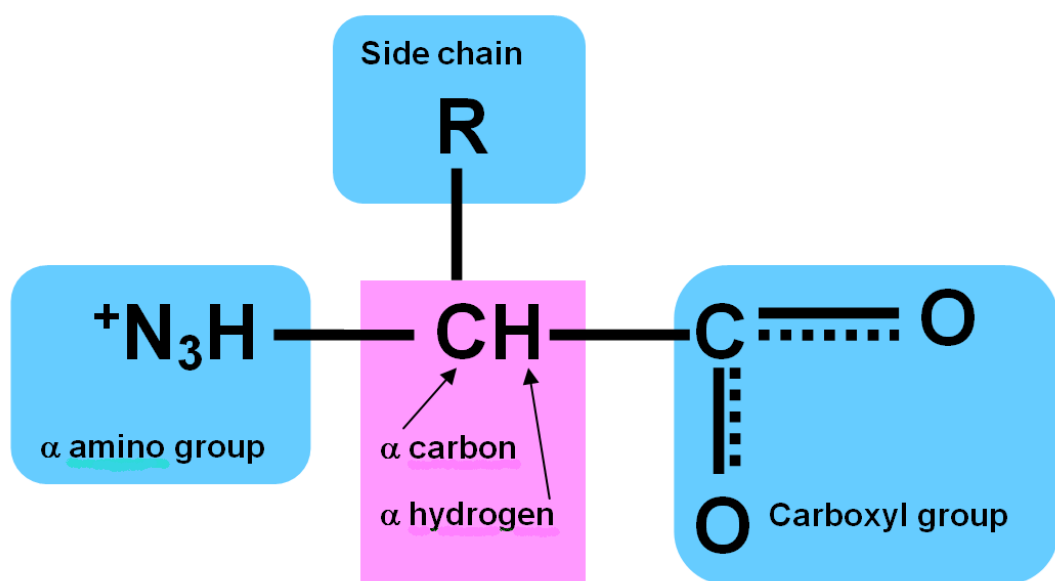


Figure 1.1. Basic structure of an amino acid.

N-C $\alpha$  bond length is 1.46 Å except for glycine and proline, which is 1.45 Å and 1.47 Å respectively (Engh and Huber, 1991). The average value of the bond lengths is shown in Figure 1.2.

The dihedral angle is also important in the protein conformation analysis. The dihedral angles for the backbone of the peptide are assigned as  $\phi$  and  $\psi$  angle for the N - C $\alpha$  bond and C $\alpha$  - C bond accordingly (Figure 1.3). Both parameters are used for the protein structure validation especially in the Ramachandran plot analysis as described by Ramachandran *et al.*, (1963) (Figure 1.4). Amino acids with torsion angles in the range  $-180 < \phi < 0^\circ$ ,  $-100 < \psi < 45^\circ$  are considered to be in the  $\alpha$ -helical region. Amino acids with torsion angles in the range  $-180 < \phi < -45^\circ$ ,  $45^\circ < \psi < -45^\circ$  are considered as  $\beta$ -sheet region. The turn region are in the range  $0^\circ < \phi < 180^\circ$  and  $-90^\circ < \psi < 90^\circ$ .

The analysis of 237,384 amino acids in 1,042 protein subunit from the PDB by Hovmoller *et al.*, (2002) found that the average values of two torsion angles,  $\phi$  and  $\psi$  of amino acids in the  $\alpha$ -helix region is within  $\pm 2^\circ$  of  $-63.8^\circ$  ( $\phi$ ) and  $-41.1^\circ$  ( $\psi$ ). The average torsion angle is  $\phi = -122^\circ$  and  $\psi = +136^\circ$  for antiparallel  $\beta$ -sheets and  $\phi = -116^\circ$ ,  $\psi = +128^\circ$  for parallel  $\beta$ -sheets (Hovmoller *et al.*, 2002).

The properties of the R group will give the unique physical and chemical properties to each amino acid. The R side chain structure and their three-letter and one-letter code are shown in Table 1.1. Amino acids can be classified according to their physical, chemical and the structure properties. One such classification was reported by Taylor (1986), Figure 1.5.



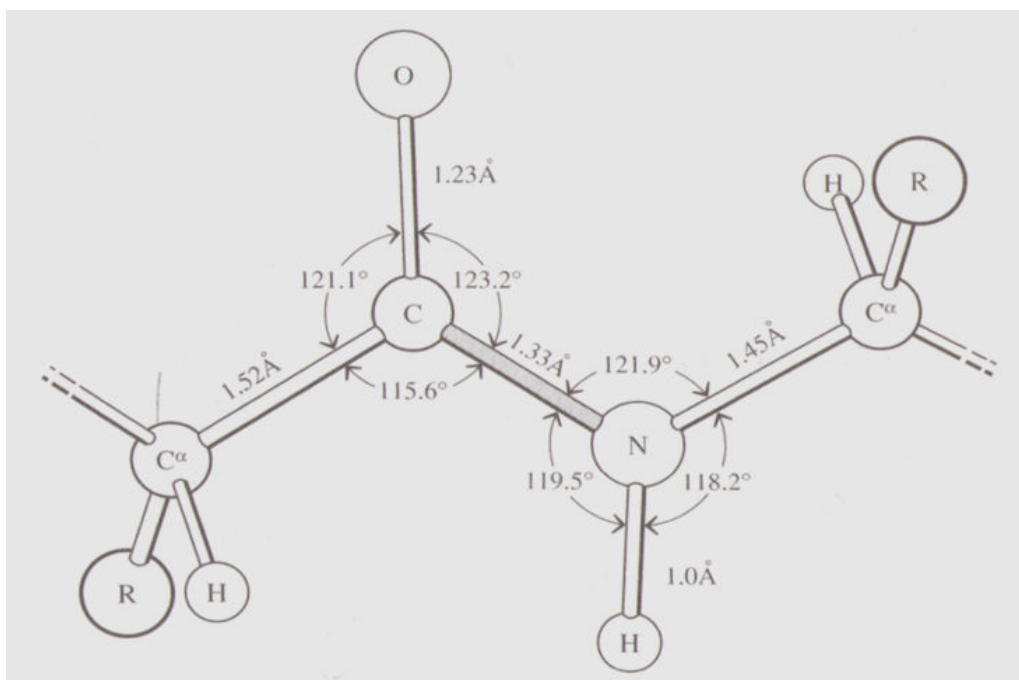


Figure 1.2. Geometry of the peptide backbone. The dimensions given are the averages observed crystallographically in amino acids and small peptides (Creighton, 2002).

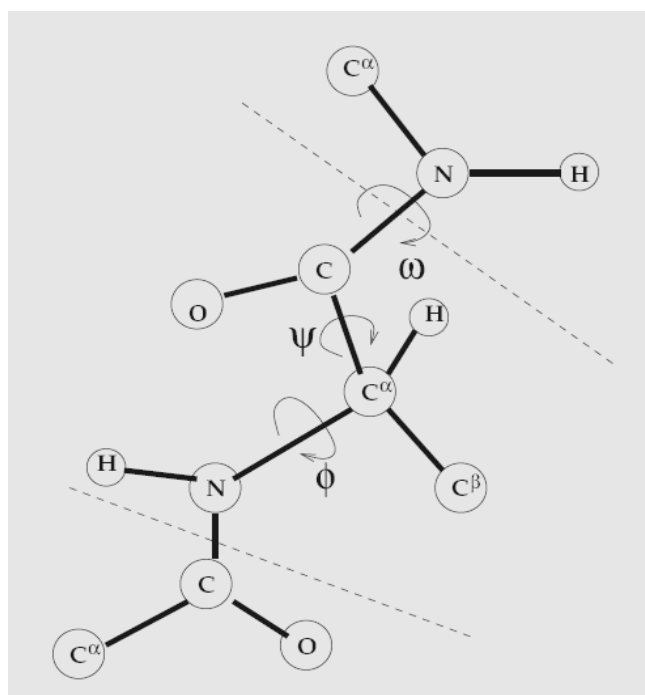


Figure 1.3. Ramachandran plot angles. The  $\phi$  and  $\psi$  angle used in the peptide dihedral angle analysis.

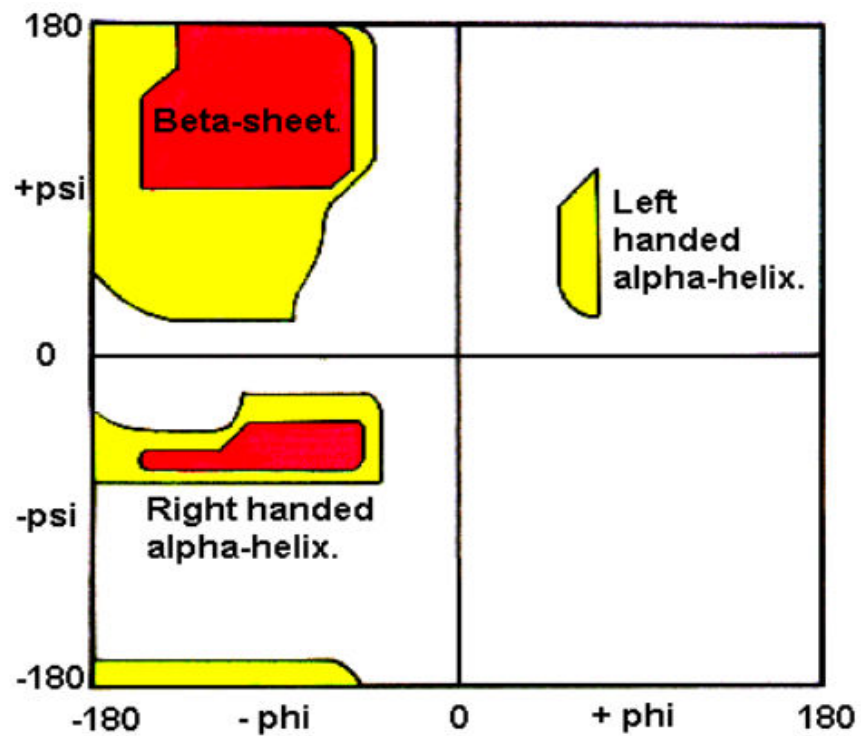


Figure 1.4 The Ramachandran plot showing the dihedral angle at N-C $\alpha$ -C, i.e. the  $\phi$  and  $\psi$  angle for the peptide backbone in the protein structure. The red color represents fully allowed region and the yellow color represents the outer limit of the  $\phi$  and  $\psi$  angle.

Table 1.1. Standard amino acid abbreviation and side chain properties

AMINO ACID	3-LETTER	1-LETTER	SIDE CHAIN
Alanine	Ala	A	CH <sub>3</sub> -
Arginine	Arg	R	HN=C(NH <sub>2</sub> )-NH-(CH <sub>2</sub> ) <sub>3</sub> -
Asparagine	Asn	N	H <sub>2</sub> N-CO-CH <sub>2</sub> -
Aspartic acid	Asp	D	HOOC-CH <sub>2</sub> -
Cysteine	Cys	C	HS-CH <sub>2</sub> -
Glutamine	Gln	Q	H <sub>2</sub> N-CO-(CH <sub>2</sub> ) <sub>2</sub> -
Glutamic acid	Glu	E	HOOC-(CH <sub>2</sub> ) <sub>2</sub> -
Glycine	Gly	G	H-
Histidine	His	H	N=CH-NH-CH=C-CH <sub>2</sub> -  _____
Isoleucine	Ile	I	CH <sub>3</sub> -CH <sub>2</sub> -CH(CH <sub>3</sub> )-
Leucine	Leu	L	(CH <sub>3</sub> ) <sub>2</sub> -CH-CH <sub>2</sub> -
Lysine	Lys	K	H <sub>2</sub> N-(CH <sub>2</sub> ) <sub>4</sub> -
Methionine	Met	M	CH <sub>3</sub> -S-(CH <sub>2</sub> ) <sub>2</sub> -
Phenylalanine	Phe	F	Phenyl-CH <sub>2</sub> -
Proline	Pro	P	-N-(CH <sub>2</sub> ) <sub>3</sub> -CH  _____
Serine	Ser	S	HO-CH <sub>2</sub> -
Threonine	Thr	T	CH <sub>3</sub> -CH(OH)-
Tryptophan	Trp	W	Phenyl-NH-CH=C-CH <sub>2</sub> -  _____
Tyrosine	Tyr	Y	4-OH-Phenyl-CH <sub>2</sub> -
Valine	Val	V	CH <sub>3</sub> -CH(CH <sub>3</sub> )-

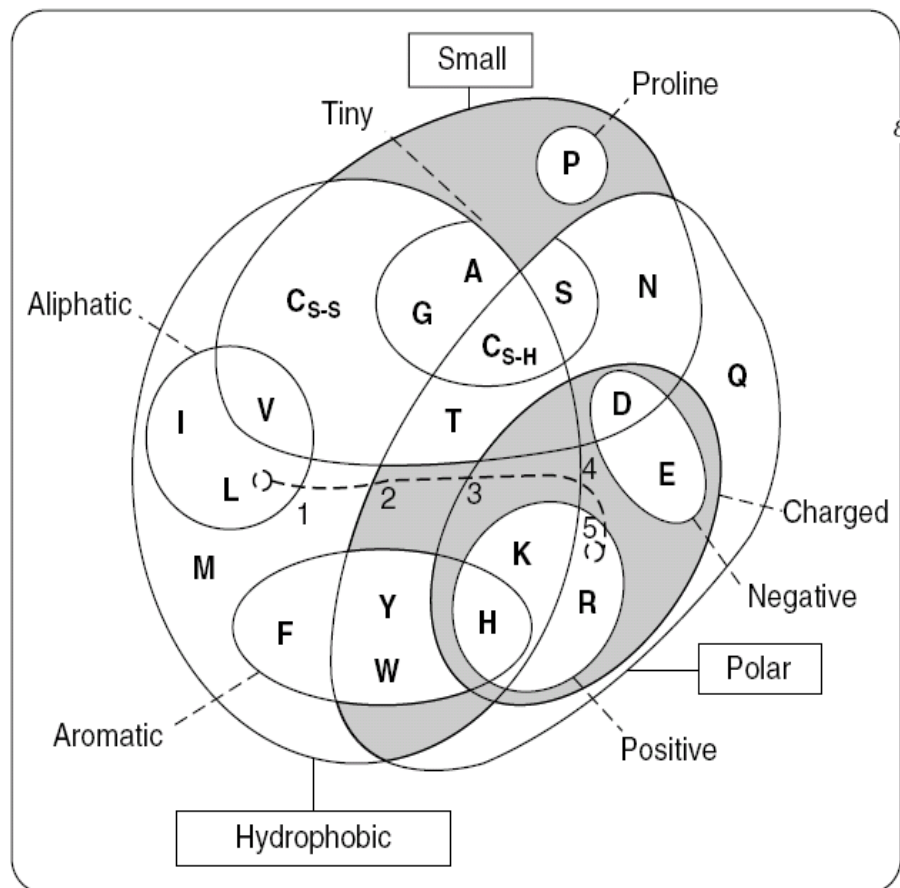


Figure 1.5. A Venn diagram illustrating the properties of amino acids (Taylor, 1986, Betts and Russell, 2003)

Most commonly, amino acids are classified according to their compatibility in aqueous environment – those that are compatible are classified as hydrophilic and those that do not, are classified as hydrophobic. The hydrophobic amino acids can be divided according to their aliphatic or aromatic side chains. Aliphatic side chains are non-reactive and are rarely involved in protein function. The aliphatic side chain amino acids are alanine, isoleucine, leucine, proline, and valine. A side chain is classified as aromatic when it contains an aromatic ring system, which is planar and the electrons are shared over the whole ring structure. The aromatic amino acids are phenylalanine, tryptophan, tyrosine and histidine.

The amino acids that can be surrounded by water have polar side chains and those buried in the protein structure usually form hydrogen bonds with other amino acids, either with other side chain or the main chain. Examples of of this group of amino acids are aspartic acid and glutamic acid are negatively charged. Lysine and arginine are positively charged. Histidine, asparagines, serine, threonine and tyrosine are all neutral.

A group of amino acid is subcategorized as small and tiny, grouped together according to their small side chain (Betts and Russell, 2003). These are alanine, cysteine, glycine, proline, serine and threonine.

The amino acid side chain properties are also used in the qualitative description of each position in a multiple alignment and the information could also be used for protein engineering in identifying the alternative amino acid that could be used as a replacement at any position for specific application. As an example, the substitution of a small side chain for a large side chain can be disastrous to the protein structure (Betts and Russell, 2003).

Betts and Russell (2003) also reported that the polar side chain such as tyrosine, histidine, arginine, lysine, aspartate, glutamate, asparagines, glutamine, serine, threonine and cystine residues are frequently found in the protein active or binding site. This shows that each amino acid, with varying side chain property plays an important role in the protein structure as well as in enzymatic reaction mechanism.

#### **1.3.1.1 Primary structure**

The primary structure is the linear order of amino acid residues along the polypeptide chain. The heteropolymer chain may contain any number and combinations of the 20 amino acids. The linear amino acid sequence is formed by covalent linkages of individual amino acids via peptide bonds. Covalent bonding between two amino acids is formed by the carbon atom from the carbonyl group of one amino acid sharing electrons with the nitrogen atom from the amino group of a second amino acid and a molecule of water is formed during this process (Alberts *et al.*, 1994). The diverse number of amino acid combinations in amino acid sequences is the basis of protein structure and function diversity.

Amino acid sequence in the protein molecule plays an important role in protein folding. As we know, the 20 amino acid side chains vary in their physico-chemical properties. Grigoriev and Kim (1999) studied the amino acid properties to identify the fold similarity between two proteins which had 15% amino acid sequence identity. The physical properties of neighboring amino acid residues in sequence at structurally equivalent position of two proteins of similar fold are often correlated even though the amino acid sequences are different (Grigoriev and Kim, 1999). Current bioinformatic tools can predict the protein secondary and tertiary

structure when the amino acid sequences of the proteins are available, because the pattern of amino acid sequence will give specific structure according to their side chain properties (Hung and Samudrala, 2003).

Organisms that live in the different temperature optima (mesophiles and thermophiles) differ in their amino acid composition and have their own protein repertoire (Table 1.2). This is to maintain the protein structure in their specific environment (Jaenicke, 2000). Analyzing amino acid sequences, enables us to understand the protein in its natural environment

### **1.3.1.2 Secondary structure**

The secondary structure is the local folding of the polypeptide chain. The spatial relationship of amino acid residues that are closed together in the primary sequence tend to form specific secondary structures.

In 1951, the secondary structure of protein was clearly understood with the availability of the X-ray protein crystallography data. The protein secondary structure is comprised of helices, beta sheets and extended structures (Pauling and Corey, 1951, Pauling *et al.*, 1951). The interactions of the side chains and main chains result in the formation of secondary structures such as alpha helices (Figure 1.6), beta sheets (Figure 1.7) and turns of the protein.

The  $\alpha$ -helix structure is formed by the hydrogen bond between the oxygen of the carbonyl group (C=O) of one turn and the hydrogen of the amide group (N-H) in the neighboring turn. One turn is formed by 3.6 residues. The distance between two turns is 0.54 nm. The helix structure can be right-handed or left-handed.



Table 1.2. Relative amino acid compositions of mesophiles and thermophiles One-letter abbreviations of amino-acid residues in brackets

	Mesophiles	Thermophiles
Charged residues (DEKRH)	24.11%	29.84%
Polar:uncharged residues (GSTNQYC)	31.15%	26.79%
Hydrophobic residues (LMIVWPAF)	44.74%	43.36%

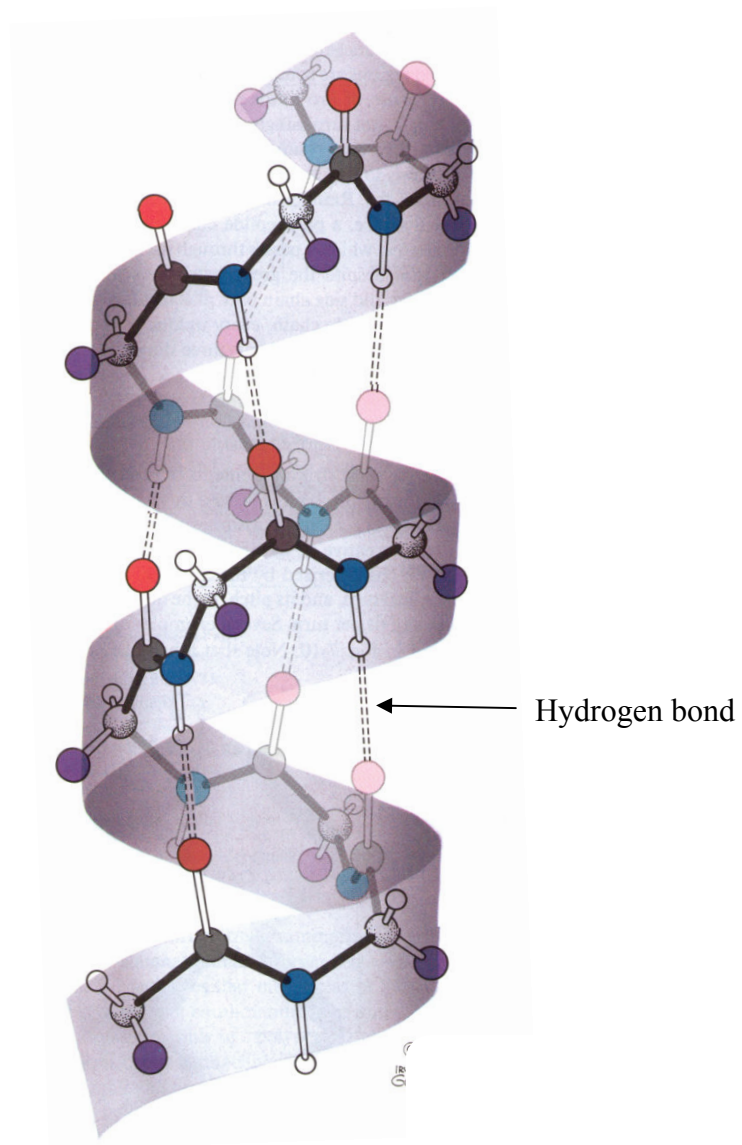


Figure 1.6. A right-handed alpha helix. Hydrogen bonds between the N – H groups and the C=O groups that are four residues back along the polypeptide chain are indicated by dashed lines (Voet and Voet, 1995, Orengo *et al.*, 1997).

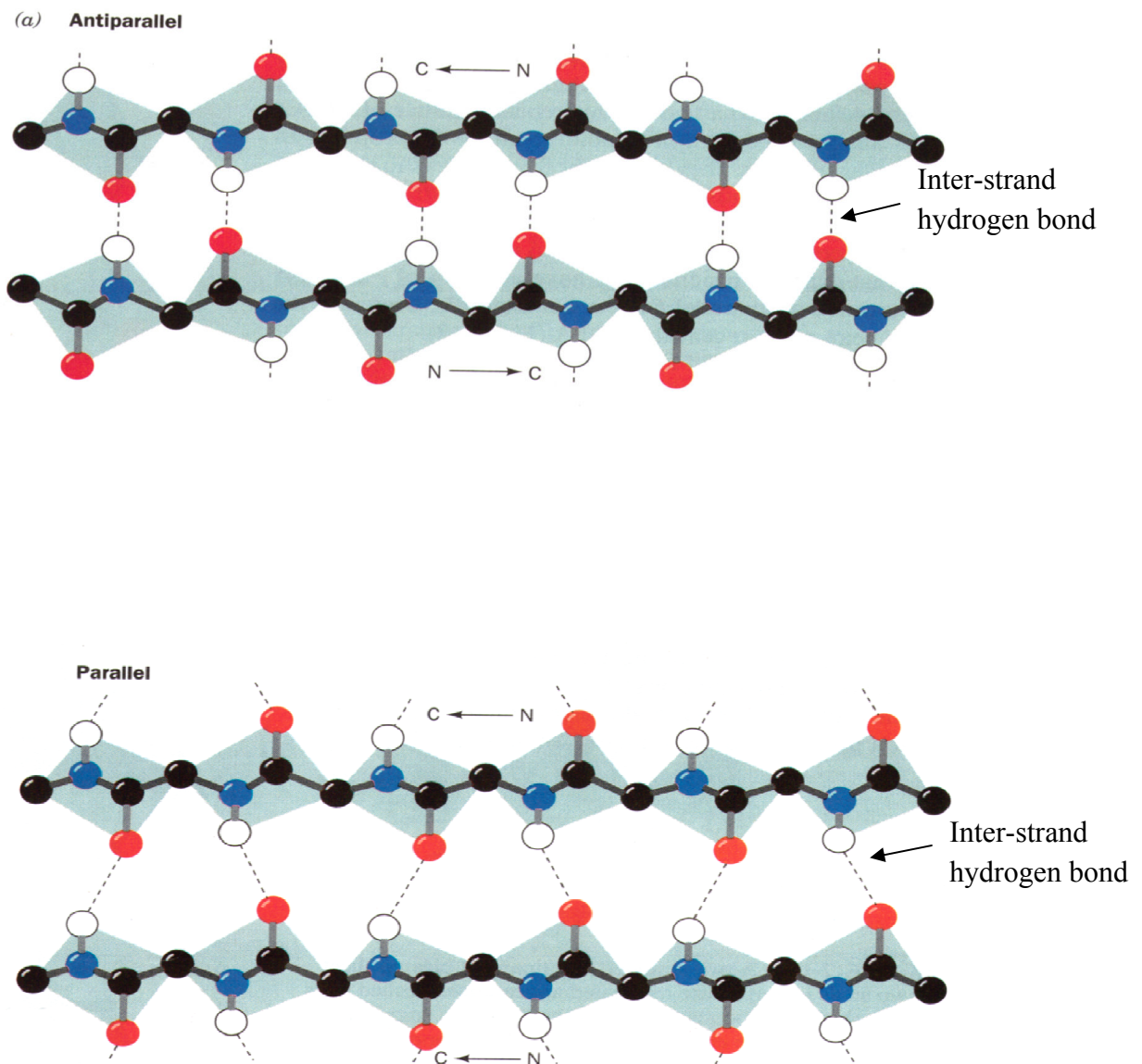


Figure 1.7. Two forms of  $\beta$ -sheets are commonly found in secondary structure, namely antiparallel and parallel  $\beta$ -sheets. The hydrogen bonding between the oxygen and hydrogen atoms at the carbonyl group and amino group are shown as dotted lines (Voet and Voet, 1995).

The right handed  $\alpha$ -helix structure is the most identifiable unit of the secondary structure. The ideal value of the dihedral angles for  $\alpha$ -helix structures is  $-57^\circ$  and  $-47^\circ$  of  $\phi$  and  $\psi$  angles respectively.

In globular proteins, helices are stable in short structure. The average length is about 12 residues and the most probable helix length is less than 6 residues. Helix structure is stable in globular protein and remains 100% helical up to temperature near the denaturation point (Dill, 1990, Kabsch and Sander, 1983).

Beta strand is the second unit of secondary structures identified by Pauling and Corey (1951). It exists as an extended conformation. Two or more  $\beta$  strands that are located close to one another tend to form additional hydrogen bonds, which stabilizes the structure. A single  $\beta$  strand is not stable. Hydrogen bonds are formed between the oxygen of a carbonyl group in amino acid of one strand and a nitrogen in the backbone of a second amino acid of another strand (Figure 1.7).

Beta-sheets can be divided into two types, the parallel or anti-parallel. A parallel beta-sheet is formed when the amino acid terminal residue of each strand points in the same direction and the anti-parallel sheet is formed when the amino termini are pointing in opposite directions. The parallel  $\beta$ -strand has dihedral angles of  $-139^\circ$  ( $\phi$ ) and  $+135^\circ$  ( $\psi$ ) respectively, while the anti-parallel  $\beta$ -strand has dihedral angles of  $-119^\circ$  ( $\phi$ ) and  $+113^\circ$  ( $\psi$ ) respectively (Whitford, 2005).

Turns are one of the elements of secondary structure. Turn conformations play a very important role in protein structure as it enables the polypeptide to change directions in the protein structure. Residues with small side chains such as glycine, aspartate, asparagines, serine, cysteine and proline are dominant in this structure (Whitford, 2005). These helices and sheets are then assembled into a compact structure recognized as tertiary structure.

### **1.3.1.3 Tertiary structure of proteins**

The tertiary protein structure is the final specific geometric shape of the protein which represents the folds of the polypeptide chain. Different fragments of secondary structure in the same chain become bonded together. Levitt and Chothia (1976) were classified globular proteins into four groups according to their secondary structure elements. The four protein classes are all  $\alpha$ -helix, all  $\beta$ -sheet,  $\alpha/\beta$ -fold and  $\alpha+\beta$  fold protein. The all  $\alpha$  and all  $\beta$  globular protein dominantly contain either  $\alpha$ -helices and  $\beta$ -sheets in their secondary structures. The  $\alpha/\beta$  proteins containing both  $\alpha$ -helices and  $\beta$ -sheets which mostly show  $\beta\alpha\beta$  unit repeats which consist of two adjacent  $\beta$  strands connected by a single  $\alpha$  helix (Russell, 2000). The  $\alpha+\beta$  globular proteins exists as  $\alpha$ -helix and  $\beta$ -strand secondary structure segments that do not mix. The  $\alpha$ -helix and  $\beta$ -strand are segregated in the polypeptide chain, a mixture of all  $\alpha$  and all  $\beta$ -sheets regions within the same polypeptide chain (Michie *et al.*, 1996).

The driving forces responsible for protein folding comes from two types of energy. The first is the bonded interaction, which includes the sum of bond, bond angle, dihedral angle and torsional angle. The second is nonbonded interaction such

as van der waals and electrostatic interaction between all pairs of atom (Oostenbrink *et al.*, 2004, Jaenicke, 2000, Dill, 1990). The bonded and nonbonded energy are all considered when using force field algorithm for computational protein structure prediction.

In protein folding, the secondary structure element and the hydrophobic interaction, which is classified as nonbonded interaction in force field algorithm is hypothesized as the major driving force in the formation of tertiary structure. The secondary structure forms first, and the chain segment then pack together by hydrophobic force (Dill, 1990) to give the final tertiary structure (Levitt, 1976). The other force fields (bonded and nonbonded interaction) also contribute in tertiary structure stability.

#### **1.3.1.4 Quarternary structure**

The quarternary structure is formed when several protein molecules interact together. A multi-subunit protein maybe composed of two or more identical polypeptides or different polypeptides. The tertiary structure aggregates to form homo or hetero-multimers (Henrick and Thornton, 1998). Quarternary structures are stabilized by interactions between residues exposed on the polypeptide surfaces within a complex. These interaction include disulfide bridges, hydrophobic interactions, charge-pair interaction and hydrogen bonding (Whitford, 2005).

#### **1.3.2 Protein structure determination**

Currently, more than 55,000 sets of atomic coordinates of protein structure are deposited in the Protein Data Bank (PDB) (<http://www.rcsb.org/pdb>) (September 2009) and the numbers are increasing daily. The protein structures were mostly

determined by X-ray protein crystallography and some of them were determined by nuclear magnetic resonance (NMR) spectroscopy.

Even though the number of protein structure deposited in the PDB keep increasing, many protein structures cannot be solved experimentally using X-ray crystallography and NMR, due to the limitations of each method. In such cases, some protein structures can be predicted from amino acid sequences using computational methods. This domain is also known as Structural Bioinformatics.

### **1.3.3 Protein structure database**

The current single worldwide depository for experimentally solved three dimensional structure of biological macromolecules is the Protein Data Bank (PDB) (Berman *et al.*, 2000). The PDB was established in Brookhaven National Laboratories in 1971. As many as 55,285 protein structures were deposited in PDB as of September 2009.

Information on protein structures can also be found in the Universal Protein Resources (UniProt), The Structural Classification of Protein (SCOP) and a semi-automatic hierarchical domain classification of protein structures in PDB (CATH) (Orengo *et al.*, 1997).

The UniProt contains information on proteins sequences and their function, and is created by combining information from Swiss-Prot and Protein Information Resources (PIR) databases (Jain *et al.*, 2009, Barker *et al.*, 1999). SCOP shows the structure and evolutionary relationship between all protein structures by manual identification (Murzin *et al.*, 1995). Protein structures with high structural similarity are classified in the same 'fold'. The term 'Superfamily' is used for proteins which probably have common evolutionary origins, whereas the term 'Family' is used for proteins that have clear evolutionary relationships.

CATH is a database of manually curated classification of protein domain structures. Protein structures are classified using a combination of automated and manual procedures. The four major levels in this classification are class, architecture, topology and homologous superfamily (Cuff *et al.*, 2009, Orengo *et al.*, 1997).

#### **1.3.4 Protein structure prediction**

The fundamentals of structural bioinformatics involved the prediction of three-dimensional protein structure from its primary sequences. The main methods currently applied to predict protein structures are homology modeling, fold recognition and *ab initio* prediction (Jones, 2000).

Zhang (2008) divided protein structure prediction to two main methods. The first is a template-based structure prediction, which requires a protein structure that has at least weak homologies to the protein of interest. The second method is a free modeling method, which does not provide appropriate templates for the prediction. In this case, the protein model is built from scratch and is called as *ab initio* or *de novo* protein structure prediction.

The similarity of the query sequence to an experimentally determined protein structure can be used as an indicator in making a decision to decide which method is to be used to predict the structure of an unknown protein. The sequence identity between the target and the template can be categorized into three classes. The first is the high homology category where the two protein sequences have 50% or above sequence identity. In the second category, the sequence identity is between 20 to 30%. This range is also known as the ‘twilight zone’ of sequence identity. The third is the ‘midnight zone’ for low sequence identity, covering the range of 8 to 10% (Rost, 1999, Rost, 1997, Chung and Subbiah, 1996).



The direct homology modeling can be applied if the target and its template amino acid sequence is 50% or greater because both structures normally share folds in the tertiary structure with a variation of 1Å in its backbone RMSD value (Gerstein and Levitt, 1998). If the amino acid sequence identity between the target and its template is below 30%, the fold recognition or threading-method can be used to assign the correct fold to the target sequence (Jones *et al.*, 1992). For proteins with a low amino acid sequence identity when compared to available determined structures, the folding cannot be assigned by threading method. In these cases, the *ab initio* method can be applied in order to assign a new fold of protein (Hardin *et al.*, 2002, Bonneau and Baker, 2001, Bonneau *et al.*, 2002).

#### **1.3.4.1 Homology modeling**

Homology modeling is the prediction of an unknown protein structure and mapping it to a known structure, usually where both proteins have evolutionarily relationship and show high amino acid sequence similarity. Proteins with high sequence similarity share similar structures. This method is usually applied to proteins which have 50% or above amino acid sequence identity to proteins with solved three dimensional (3-D) structures in the PDB.

Once a suitable template for homology modeling is identified, Modeller is used to construct three-dimensional models (Sali and Blundell, 1993, Sali and Overington, 1994, Eswar *et al.*, 2006, Eswar *et al.*, 2007, Marti-Renom *et al.*, 2000). Modeller performs comparative modeling by predicting the 3-D structure of a given protein sequence (target) based primarily on its alignment to one or more known structures of proteins (templates). The prediction process involves fold assignment, target-template alignment, model building and model evaluation.

#### 1.3.4.2 Fold Recognition

The fold-recognition approach in protein structure prediction involves the identification of experimentally determined protein structures that can accommodate the target protein sequence. The structure based fold recognition is also known as the threading method. Threading works by computing a scoring function that performs the best fit of a sequence against an experimentally determined structure. Typically, a threading computer program comprises of four components. First is the template structure, which is usually a protein structure deposited in the Protein Data Bank database. Second is the evaluation of the compatibility between the target sequence and the template fold. This is followed by the algorithm to compute the optimal alignment between the target sequence and the template structure. Finally the statistical significance is estimated based on the ranking (Fischer *et al.*, 1996).

The target-template compatibility functions typically using two types of structural feature, the pattern of secondary structure element and local environment classes. The local environment classes are the combination of solvent accessibility, polarity of the side chain environment and local backbone conformation (Rost *et al.*, 1997).

Cymerman *et al.*, (2004) defined the fold-recognition algorithms to two main types namely the sequence based fold recognition and structure based fold recognition method. The principle of the sequence based fold recognition method is to search structurally characterized proteins from the protein structure database that exhibited significant sequence similarity to the target protein. This step is done by using BLAST or PSI-BLAST in order to identify suitable template from the protein data bank.